

NOVA: Sparse Control, Dense Synthesis for Pair-Free Video Editing

Supplementary Material

A. Details of Experiments on Naive Multi-Keyframe Guidance

To explore the feasibility of multi-keyframe guidance for temporally consistent video editing, we conduct preliminary experiments using the reference-to-video diffusion model WAN VACE. VACE accepts two auxiliary inputs: (1) a *reference video*, which provides semantic or structural guidance (e.g., edited RGB frames, depth maps, or pose skeletons), and (2) a *mask video*, which indicates per-frame regions to be preserved (black pixels) or regenerated (white pixels).

Our initial attempt adopts a straightforward multi-keyframe strategy: we construct the reference video by placing the user-edited keyframes at their corresponding timestamps, while filling all non-keyframe positions with a uniform gray frame (RGB value 127). Concurrently, the mask video is set to black at keyframes (instructing the model to faithfully reconstruct the edited content) and white elsewhere (prompting content generation based on the reference). However, this naive setup leads to severe temporal flickering and inconsistent motion in both the 1.3B and 14B parameter variants of VACE. We hypothesize that this instability arises because VACE was not trained on such sparse, discontinuous reference signals, resulting in hallucinated textures and incoherent motion interpolation between keyframes.

To mitigate this issue, we design a more stable reference configuration. Specifically, we retain the edited RGB image only at the first keyframe (typically the first frame of the video) in the reference video. For all subsequent keyframes, we replace the edited RGB content with the corresponding *depth maps* derived from the original edited frames. Depth maps are chosen because they encode structural and motion cues with significantly less high-frequency appearance variation than RGB images, thereby reducing the risk of introducing inconsistent visual signals. Meanwhile, the mask video is simplified: only the first frame is set to black (to anchor the edited appearance), and all remaining frames (including other keyframes—are set to white), indicating that the model should synthesize these frames using the reference guidance.

This refined setup effectively decouples appearance and motion guidance: the model learns the target appearance from the first edited frame, while leveraging depth-derived structural cues from later keyframes to preserve the original video’s motion dynamics. Under this configuration, VACE produces visibly more stable outputs with reduced flickering, enabling us to meaningfully analyze the limitations of

sparse multi-keyframe editing

B. Details of Training Pipelines

Algorithm 1 Source Fidelity Pipeline (Pseudo-Source Generation)

Require: Target video $\mathcal{X} = \{\mathbf{x}_t\}_{t=0}^T$, Filler video pool \mathcal{P}
Ensure: Pseudo-source video $\tilde{\mathcal{X}} = \{\tilde{\mathbf{x}}_t\}_{t=0}^T$

- 1: Sample filler video $\mathcal{Y} = \{\mathbf{y}_t\} \sim \mathcal{P}$
- 2: Init mask shape $S \in \{\text{rect, ellipse, poly}\}$
- 3: Init position \mathbf{p}_0 , motion \mathbf{v} , rotation θ
- 4: **for** $t = 0$ to T **do**
- 5: Update mask pose: $\mathbf{p}_t \leftarrow \mathbf{p}_{t-1} + \mathbf{v}$, $\theta_t \leftarrow \theta_{t-1} + \Delta\theta$
- 6: Handle boundary collisions and bounce \mathbf{v} if needed
- 7: Generate binary mask \mathbf{m}_t from $S, \mathbf{p}_t, \theta_t$
- 8: Extract patch from \mathbf{y}_t (using ping-pong loop)
- 9: $\tilde{\mathbf{x}}_t \leftarrow \mathbf{m}_t \odot \mathbf{y}_t + (\mathbf{1} - \mathbf{m}_t) \odot \mathbf{x}_t$
- 10: **end for**
- 11: **return** $\tilde{\mathcal{X}}$

Algorithm 2 Anchored Control Pipeline (Degraded Reference)

Require: Target video \mathcal{X} , Indices set $\{0, \dots, T\}$
Ensure: Degraded reference video $\hat{\mathcal{X}}$

- 1: $\mathcal{K} \leftarrow \{0, T\}$ ▷ Always include start/end
- 2: Sample N random indices from $(0, T)$ into \mathcal{K}
- 3: Sort $\mathcal{K} = \{k_0, k_1, \dots, k_{|\mathcal{K}|}\}$
- 4: **for** each keyframe index $k \in \mathcal{K}$ **do**
- 5: $\hat{\mathbf{x}}_k \leftarrow \mathbf{x}_k$
- 6: **if** $\text{random}() < 0.5$ **then** ▷ Geometric Degradation
- 7: Apply random Zoom-Stretch affine transform to $\hat{\mathbf{x}}_k$
- 8: **end if**
- 9: **if** $\text{random}() < 0.5$ **then** ▷ Appearance Degradation
- 10: Define local region \mathbf{b}_k (random blob)
- 11: $\hat{\mathbf{x}}_k \leftarrow (\mathbf{1} - \mathbf{b}_k) \odot \hat{\mathbf{x}}_k + \mathbf{b}_k \odot \text{Blur}(\hat{\mathbf{x}}_k)$
- 12: **end if**
- 13: **end for**
- 14: **for** $t = 0$ to T **do** ▷ Temporal Interpolation
- 15: Find neighbors $k_i, k_{i+1} \in \mathcal{K}$ such that $k_i \leq t \leq k_{i+1}$
- 16: $\alpha \leftarrow (t - k_i) / (k_{i+1} - k_i)$
- 17: $\hat{\mathbf{x}}_t \leftarrow (1 - \alpha)\hat{\mathbf{x}}_{k_i} + \alpha\hat{\mathbf{x}}_{k_{i+1}}$
- 18: **end for**
- 19: **return** $\hat{\mathcal{X}}$

Add a wooden house near the path

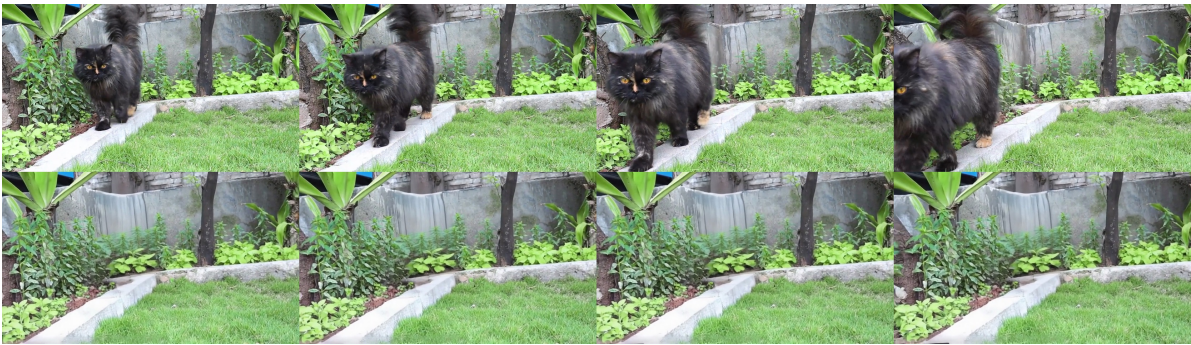


Add a window on the wall



Figure 12. More visual results (1).

Remove the black cat



Remove the man and the motor



Figure 13. More visual results (2).

Remove the bird on the beach



Add a tower by the highway



Figure 14. More visual results (3).

Add sand drawing "hello" on the beach



Remove the mountains



Figure 15. More visual results (4).